

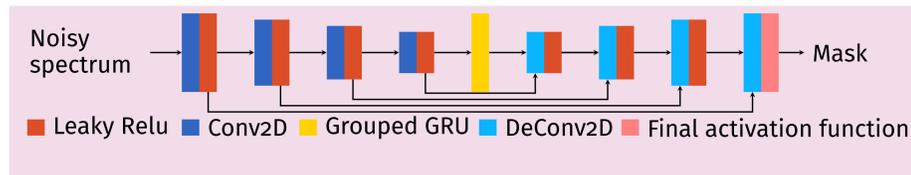
PHASE RECONSTRUCTION IN SINGLE CHANNEL SPEECH ENHANCEMENT BASED ON PHASE GRADIENTS AND ESTIMATED CLEAN-SPEECH AMPLITUDES

Yanjue Song, Nilesh Madhu

IDLab, Department of Electronics and Information Systems, Ghent University - imec, Belgium

{yanjue.song, nilesh.madhu}@ugent.be

1 Introduction



$$S(l, m) = |S(l, m)| \exp(j\Phi_S(l, m))$$

- Phase estimation in single-channel speech enhancement
- Taking the noisy phase: the MMSE-sense optimal
- Complex mask/mapping
- Challenge:** no pattern observed
- Pattern exists in phase gradients

$$\Delta_t \Phi(l, m) = \Phi(l, m) - \Phi(l, m-1) \quad (1)$$

$$\Delta_f \Phi(l, m) = \Phi(l, m) - \Phi(l-1, m) \quad (2)$$

→ Estimate phase from clean magnitude (phase retrieval)

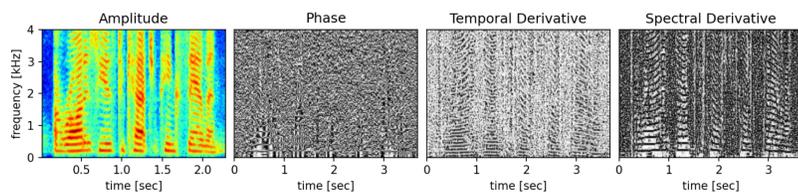
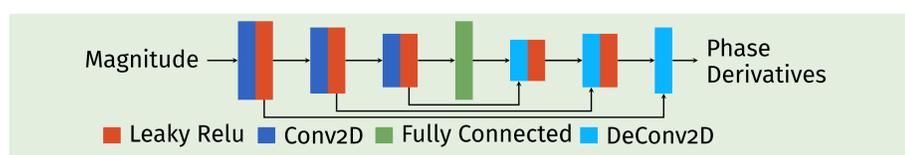


Figure 1: The amplitude, phase, and phase derivative of the clean speech.

2 Motivation

- Can we utilise the gradient information in speech enhancement?
 - Phase retrieval solution: too artificial
 - What's missing: the initial phase estimate
 - Idea:** ground it by the initial phase estimate
 - Our method: fuse
 - (a) Phase estimate from $\Delta_t \Phi(l, m)$ (temporal derivative)
 - (b) Phase estimate from $\Delta_f \Phi(l, m)$ (spectral derivative)
 - (c) **Initial phase estimate**
- to get *one consistent phase estimate* for the enhanced speech

3 Phase Derivative Estimation



- Loss function: $\mathcal{L}_* = \sum_{l,m} (1 - \cos(\Delta_* \widehat{\Phi}(l, m) - \Delta_* \Phi(l, m)))$, $\star \in \{t, f\}$
- Training scheme: **matched** or **agnostic**?

4 Phase Reconstruction

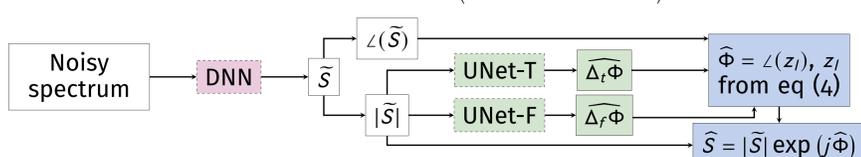
- Goal: one consistent phase estimate based on the three different estimates
- Minimise the distance between this final estimate z_l and all sources
- Translated to the cost function:

$$\mathcal{J}(z_l) = \underbrace{\|z_l - \widehat{V}_l \circ \widehat{S}_{l-1}\|_{\Lambda_l}^2}_{\text{distance to } \widehat{\Phi} \text{ from } \Delta_t \Phi} + \underbrace{\|D_l z_l\|_{\Gamma_l}^2}_{\text{distance to } \widehat{\Phi} \text{ from } \Delta_f \Phi} + \underbrace{\|z_l - \widehat{S}_l\|_{\Omega_l}^2}_{\text{distance to } \widehat{S}} \quad (3)$$

where z is the clean speech estimate.

→ The optimal solution:

$$\widehat{z}_l = (\Lambda_l + D_l^H \Gamma_l D_l + \Omega_l)^{-1} (\Lambda_l (\widehat{V}_l \circ \widehat{S}_{l-1}) + \Omega_l \widehat{S}_l) \quad (4)$$



5 Evaluation Results

Dataset

- Training: DNS challenge 2021, 140 hours
- Test: DNS challenge 2020, synthetic test set

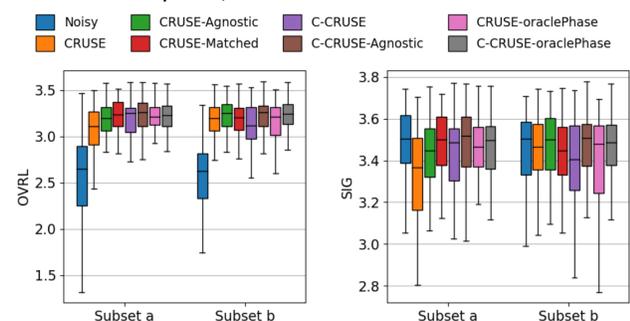
Objective metrics

Table 1: Averaged instrumental metrics on test set. Best results in bold.

Method	segSNR [dB]	STOI	DNSMOS	
			OVRL	SIG
Noisy	6.87	0.87	2.53	3.33
CRUSE	13.74	0.93	3.10	3.36
CRUSE-Agnostic	14.30	0.93	3.17	3.43
CRUSE-Matched	14.19	0.93	3.17	3.44
C-CRUSE	13.92	0.93	3.14	3.40
C-CRUSE-Agnostic	14.45	0.93	3.20	3.45
CRUSE-OraclePhase	14.51	0.94	3.17	3.43
C-CRUSE-OraclePhase	14.77	0.94	3.20	3.45

DNSMOS Distribution

- subset a) mixtures with *stationary* or *short-term stationary* noise
- subset b) mixtures with *sparse, transient* noise



- Phase enhancement
 - improves all metrics
 - is comparable to using *oracle* phase
- Boosts signal quality in poor SNR conditions
- For stationary noise: SE-Matched > SE-Agnostic
- For sparse noise: SE-Agnostic > SE-Matched

Spectrogram Samples

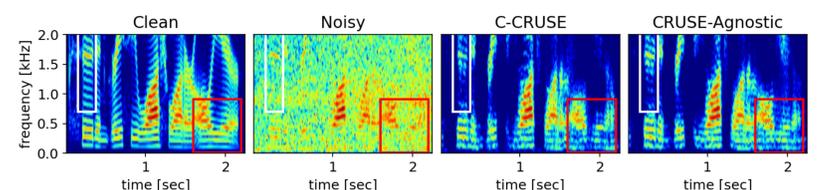


Figure 2: Noisy signal: Street noise, -2 dB.

- Similar performance in high SNR regions
- More continuous harmonics by the proposed phase reconstruction

6 Conclusions

- Incorporating initial phase estimate for natural-sounding output
- Improvement reflected in objective audio quality metrics
- Compatible with real- or complex-domain methods
- Matched/agnostic methods suit different noise types

